

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 09-160928

(43)Date of publication of application : 20.06.1997

(51)Int.Cl.

G06F 17/30

(21)Application number : 07-320544

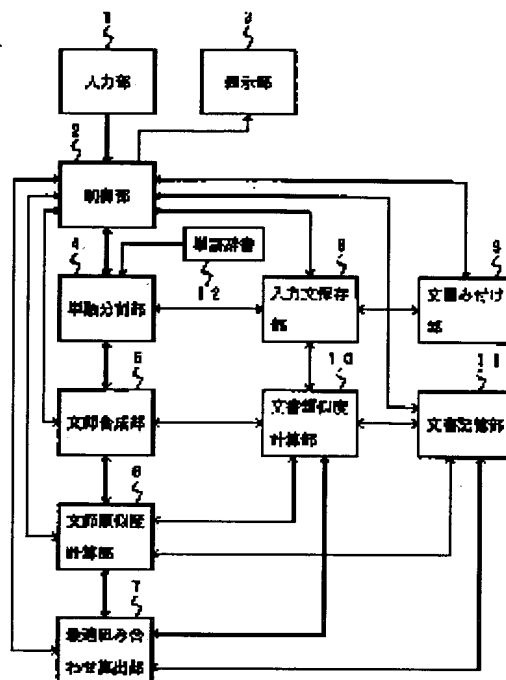
(71)Applicant : TOSHIBA CORP

(22)Date of filing : 08.12.1995

(72)Inventor : YAJIMA MASATO
KOYAMA NORIKO

(54) DOCUMENT RETRIEVING METHOD AND DEVICE THEREFOR

(57)Abstract:

PROBLEM TO BE SOLVED: To retrieve a document having contents similar to that of the document that a user desires.**SOLUTION:** This device is provided with a document similarity calculation part 10 finding the degree of similarity of the whole of the document between the input document to be the retrieval key preserved in an input sentence preservation part 8 and each document (retrieval object document) stored in a document storage part 11. When a retrieval request command is inputted from an input part 1, a sentence is taken out as a retrieval input sentence from the input sentence preservation part 8, a document is taken out as a retrieval object document from the document storage part 11 and the degree of similarity for every paragraph between the retrieval input sentence and each sentence of the retrieval object document is made to be calculated by a paragraph similarity calculation part 6. From the combination result of an optimum combination calculation part 7 based on the calculation result, the degrees of similarity of the whole of each sentence between the retrieval input sentence and each sentence of the retrieval object document are found and are successively added. By performing this operation for all the sentences in the input document, the degree of similarity of the whole of the document between the input document and the retrieval object document is found and a document retrieval is performed based on the similarity.

LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-160928

(43) 公開日 平成9年(1997)6月20日

(51) Int.Cl.⁶

G 0 6 F 17/30

識別記号

庁内整理番号

F I

G 0 6 F 15/403

技術表示箇所

3 4 0 Z

3 5 0 C

審査請求 未請求 請求項の数 6 O L (全 14 頁)

(21) 出願番号 特願平7-320544

(22) 出願日 平成7年(1995)12月8日

(71) 出願人 000003078

株式会社東芝

神奈川県川崎市幸区堀川町72番地

(72) 発明者 矢島 真人

神奈川県川崎市幸区柳町70番地 株式会社

東芝柳町工場内

(72) 発明者 小山 紀子

神奈川県川崎市幸区柳町70番地 株式会社

東芝柳町工場内

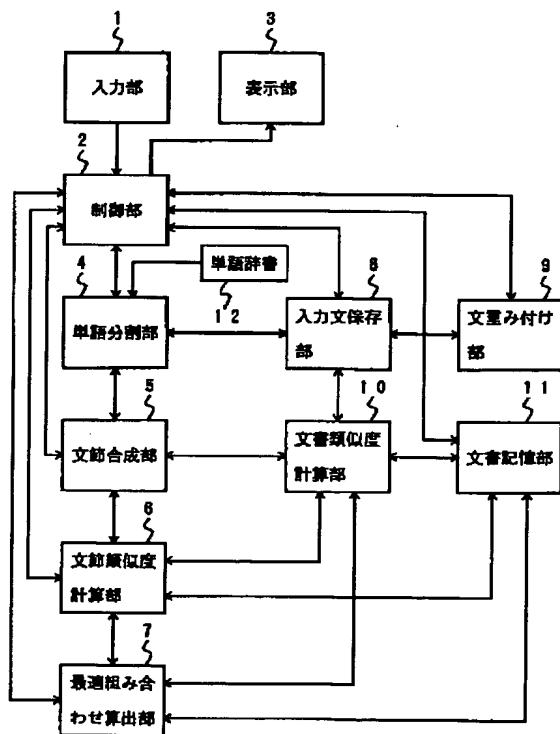
(74) 代理人 弁理士 鈴江 武彦

(54) 【発明の名称】 文書検索方法及び装置

(57) 【要約】

【課題】 ユーザが望む文書と内容が類似した文書を検索できるようにする。

【解決手段】 入力文保存部 8 に保存された検索キーとなる入力文書と、文書記憶部 11 に記憶されている各文書（検索対象文書）との間の文書全体の類似度を求める文書類似度計算部 10 を設け、入力部 1 から検索要求コマンドが入力された場合には、入力文保存部 8 から 1 文を検索入力文として取り出すと共に、文書記憶部 11 から 1 文書を検索対象文書として取り出し、検索入力文と検索対象文書の各文との間の文節ごとの類似度を文節類似度計算部 6 により計算させ、その計算結果に基づく最適組み合わせ算出部 7 の組み合わせ結果から、検索入力文と検索対象文書の各文との間のそれぞれ文全体の類似度を求めて順次加算する操作を、入力文書のすべての文について行うことで、入力文書と検索対象文書との間の文書全体の類似度を求め、その類似度をもとに文書検索を行う。



【特許請求の範囲】

【請求項1】 指定された文書中のすべての文を検索入力文として、検索対象として用意されている文書中の各文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算して、その文全体の類似度を求め、

その文全体の類似度を順次加算した値を前記指定文書と検索対象文書との間の文書全体の類似度として求める操作を、検索対象となるすべての文書について順次行うことにより、

前記求めた文書全体の類似度をもとに、前記検索対象となる複数の文書の中から前記指定文書と類似した文書を検索することを特徴とする文書検索方法。

【請求項2】 複数の文書を記憶しておくための文書記憶手段と、

検索入力文と検索対象文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算する文構成単位類似度計算手段と、

この文構成単位類似度計算手段の計算結果をもとに、前記検索入力文と検索対象文との間で最も類似している文構成単位の組み合わせを求める最適組み合わせ算出手段と、

指定された文書から1文を検索入力文として取り出すと共に、前記文書記憶手段から1文書を検索対象文書として取り出して、前記検索入力文と前記検索対象文書の各文との間の文構成単位の類似度を、前記文構成単位類似度計算手段により計算させ、その計算結果に基づく前記最適組み合わせ算出手段の組み合わせ結果から、前記検索入力文と前記検索対象文書の各文との間のそれぞれ文全体の類似度を求めて順次加算する操作を、前記指定文書のすべての文について行うことで、前記指定文書と検索対象文書との間の文書全体の類似度を求める文書類似度計算手段とを具備し、

前記指定文書との文書全体の類似度を前記文書記憶手段内の検索対象となる複数の文書のそれぞれについて求めることにより、当該文書全体の類似度をもとに、前記検索対象となる複数の文書の中から前記指定文書と類似した文書を検索することを特徴とする文書検索装置。

【請求項3】 単語辞書内の各単語ごとにその重要度を付与しておき、

指定された文書で使用されている各単語の重要度を前記単語辞書から求めることで、その各単語の重要度に応じて前記指定文書中の各文に重み付けを行い、

前記指定文書中のすべての文を検索入力文として、検索対象として用意されている文書中の各文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算して、その文全体の類似度を求め、

その文全体の類似度と当該文に付された重みとの乗算値を順次加算した値を前記指定文書と検索対象の文書との間の文書全体の類似度として求める操作を、検索対象と

なるすべての文書について行うことにより、

前記求めた文書全体の類似度をもとに、前記検索対象となる複数の文書の中から前記指定文書と類似した文書を検索することを特徴とする文書検索方法。

【請求項4】 重要度を示す情報が付された複数の単語を登録しておくための単語辞書と、

この単語辞書の示す、指定された文書で使用されている各単語の重要度に応じて前記指定文書中の各文に重み付けを行う文重み付け手段と、

複数の文書を記憶しておくための文書記憶手段と、

検索入力文と検索対象文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算する文構成単位類似度計算手段と、

この文構成単位類似度計算手段の計算結果をもとに、前記検索入力文と検索対象文との間で最も類似している文構成単位の組み合わせを求める最適組み合わせ算出手段と、

前記指定文書から1文を検索入力文として取り出すと共に、前記文書記憶手段から1文書を検索対象文書として取り出して、前記検索入力文と前記検索対象文書の各文との間の文構成単位の類似度を、前記文構成単位類似度計算手段により計算させ、その計算結果に基づく前記最適組み合わせ算出手段の組み合わせ結果から、前記検索入力文と前記検索対象文書の各文との間のそれぞれ文全体の類似度を求め、その文全体の類似度と前記文重み付け手段により当該文に付された重みとの乗算値を順次加算する操作を、前記指定文書のすべての文について行うことで、前記指定文書と検索対象文書との間の文書全体の類似度を求める文書類似度計算手段とを具備し、前記指定文書との文書全体の類似度を前記文書記憶手段内の検索対象となる複数の文書のそれぞれについて求めることにより、当該文書全体の類似度をもとに、前記検索対象となる複数の文書の中から前記指定文書と類似した文書を検索することを特徴とする文書検索装置。

【請求項5】 指定された文書で使用されている各単語の頻度を求めることで、その各単語の頻度に応じて前記指定文書中の各文に重み付けを行い、

前記指定文書中のすべての文を検索入力文として、検索対象として用意されている文書中の各文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算して、その文全体の類似度を求め、

その文全体の類似度と当該文に付された重みとの乗算値を順次加算した値を前記指定文書と検索対象の文書との間の文書全体の類似度として求める操作を、検索対象となるすべての文書について行うことにより、

前記求めた文書全体の類似度をもとに、前記検索対象となる複数の文書の中から前記指定文書と類似した文書を検索することを特徴とする文書検索方法。

【請求項6】 指定された文書で使用されている各単語の頻度を求めることで、その各単語の頻度に応じて前記

指定文書中の各文に重み付けを行う文重み付け手段と、複数の文書を記憶しておくための文書記憶手段と、検索入力文と検索対象文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算する文構成単位類似度計算手段と、

この文構成単位類似度計算手段の計算結果をもとに、前記検索入力文と検索対象文との間で最も類似している文構成単位の組み合わせを求める最適組み合わせ算出手段と、

前記指定文書から 1 文を検索入力文として取り出すと共に、前記文書記憶手段から 1 文書を検索対象文書として取り出して、前記検索入力文と前記検索対象文書の各文との間の文構成単位の類似度を、前記文構成単位類似度計算手段により計算させ、その計算結果に基づく前記最適組み合わせ算出手段の組み合わせ結果から、前記検索入力文と前記検索対象文書の各文との間のそれぞれ文全体の類似度を求め、その文全体の類似度と前記文重み付け手段により当該文に付された重みとの乗算値を順次加算する操作を、前記指定文書のすべての文について行うことで、前記指定文書と検索対象文書との間の文書全体の類似度を求める文書類似度計算手段とを具備し、前記指定文書との文書全体の類似度を前記文書記憶手段内の検索対象となる複数の文書のそれぞれについて求めることにより、当該文書全体の類似度をもとに、前記検索対象となる複数の文書の中から前記指定文書と類似した文書を検索することを特徴とする文書検索装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、検索の対象となる大量の文書中からユーザ指定の文書と類似している文書を検索するのに好適な文書検索方法及び装置に関する。

【0002】

【従来の技術】日本語ワードプロセッサや光学的文字読取装置（OCR）などの文書入力装置の普及に伴い、従来紙やマイクロフィルムという形態で保存していた文書を、電子化して磁気ディスクや光磁気ディスクなどの外部記憶装置に保存するようになってきた。また大量の電子化された文書データから、いかにユーザが要求する文書を検索するかというテキスト検索技術も開発され発展しつつある。

【0003】従来の検索技術では、ユーザが文書を検索するために検索キーとなる文字列や単語などを入力して、検索キーが含まれる文書を検索するという方式をとることが多く、検索キーである文字列や単語が検索文書内でどのように使われているかを問わないので、ユーザが望む文書以外に、意味のない文書を検索してしまうという問題があった。

【0004】これに対して、特願平 5 - 2 1 2 5 1 5 号にて提案されているように、文を一定の文構成単位、例えば文節単位に分割して文節ごとに類似度を計算し、最

も類似している文節の組み合わせを求めるという検索方式がある。この方式では、入力した文（ユーザ指定の文）に対して文全体が類似している文を検索できるという特長がある。これによると、単なる文字列や単語による検索と異なり 1 文を検索キーとするため、ユーザが期待する文章表現と類似した文を含む文書を検索することができる。

【0005】しかし、この方式では、例えばユーザが特定の文書と内容が類似した文書を検索したいと考えた場合、当該特定の文書内の各文に対して類似した文を含む文書を検索することはできるが、文書全体が類似した文書とは限らないため、やはり意味のない文書を検索してしまうという可能性を免れなかった。

【0006】

【発明が解決しようとする課題】上記したように従来の文書検索方式では、ユーザ指定の 1 文と類似した文を含む文書を検索することはできるものの、特定の文書に対して類似した文書を検索することはできなかった。

【0007】本発明は上記事情を考慮してなされたものでその目的は、特定の文書（ユーザの入力した文書あるいはユーザの指定した文書）内の各文と類似した文を含む文書を検索するだけでなく、文書全体がどれだけ類似しているかを総合的に判断することで、当該特定の文書と内容が類似した文書を検索することができる文書検索方法及び装置を提供することにある。

【0008】

【課題を解決するための手段】本発明の第 1 の観点に係る文書検索方法は、指定文書（ユーザが入力した文書あるいはユーザが指定した文書）の中のすべての文を検索入力文として、検索対象として用意されている文書中の各文との間で一定の文構成単位ごとに比較することで当該文構成単位の類似度を計算して、その文全体の類似度を求め、その文全体の類似度を順次加算した値を上記指定文書と検索対象文書との間の文書全体の類似度として求める操作を、検索対象となるすべての文書について行うことにより、この文書全体の類似度をもとに検索対象となる複数の文書の中から上記指定文書と類似した文書を検索することを特徴とする。

【0009】本発明の第 2 の観点に係る文書検索方法は、単語辞書内の各単語ごとに重要度を付与しておき、指定文書で使用されている各単語の重要度を上記単語辞書から求めることで、その各単語の重要度に応じて指定文書中の各文に重み付けを行い、上記第 1 の観点に係る文書検索方法において、文全体の類似度を順次加算した値を文書全体の類似度とするのに代えて、文全体の類似度と当該文に付された重みとの乗算値を順次加算した値を文書全体の類似度とするようにしたことを特徴とする。

【0010】本発明の第 3 の観点に係る文書検索方法は、指定文書で使用されている各単語の頻度を求めるこ

とで、その各単語の頻度に応じて前記指定文書中の各文に重み付けを行い、上記第1の観点に係る文書検索方法において、文全体の類似度を順次加算した値を文書全体の類似度とするのに代えて、文全体の類似度と当該文に付された重みとの乗算値を順次加算した値を文書全体の類似度とするようにしたことを特徴とする。

【0011】本発明においては、指定文書中のすべての文を検索入力文として、検索対象文書中の各文に対して1文全体の類似度を求めて順次加算した値を文書全体の類似度とすることで、指定文書と検索対象文書とが文書全体としてどれだけ類似しているかを、当該文書全体の類似度で表すことができるため、この文書全体の類似度をすべての検索対象文書について求めることにより、この文書全体の類似度をもとにすべての検索対象文書の中から指定文書と類似した文書を検索することが可能となる。

【0012】また本発明においては、指定文書で使用されている各単語の重要度または使用頻度に応じて当該指定文書中の各文に重み付けを行い、その各文ごとの重みをその文（検索対象文）と検索入力文との文全体の類似度に反映させることで、すべての検索対象文書の中から指定文書と類似した文書であって且つ指定文書中で使用している重要度の高い単語を多数含む、あるいは指定文書中で使用している単語の使用数の多い文書を検索することが可能となる。

【0013】このように本発明によれば、ユーザ指定の文書と文書全体が類似した文書を検索することができるために、ユーザが特定の文書と類似した文書を検索したい場合に、従来のように当該特定の文書内で使われている単語やキーワードを選び出して検索したり、重要な文を選び出して類似した文を含む文書を検索したりする手間は一切いらなくなり、単に手元にある文書をそのまま指定文書として入力したり、記憶手段内の文書群の中から所望の文書を指定文書として呼び出すだけで、ユーザが意図する文書を検索することが可能となる。

【0014】

【発明の実施の形態】以下、本発明の実施の形態につき図面を参照して説明する。図1は本発明の一実施形態に係る文書検索装置の全体構成を示すブロック図である。

【0015】図1の文書検索装置は例えばパーソナルコンピュータ等の情報処理装置により実現されるもので、入力部1、制御部2、表示部3、単語分割部4、文節合成部5、文節類似度計算部6、最適組み合わせ算出部7、入力文保存部8、文重み付け部9、文書類似度計算部10、文書記憶部11及び単語辞書12から構成される。

【0016】入力部1は、文字列や種々の制御コード（コマンド等）を入力するためのもので、例えばキーボードまたはOCR（光学的文字読取装置）からなる入力装置である。

【0017】制御部2は装置全体を司るもので、例えば中央処理ユニット（CPU）である。ここで制御部2は、文書検索処理を実行するための装置内各部（単語分割部4、文節合成部5、文節類似度計算部6、最適組み合わせ算出部7、文重み付け部9、文書類似度計算部10）の制御を行う。

【0018】表示部3は、入力文の文字列や検索対象の文書などを表示するためのもので、例えばCRTディスプレイ装置または液晶表示装置（等のフラットパネルディスプレイ装置）である。

【0019】入力文保存部8は、検索のキーとなる入力文書（指定文書）を一時的に保存するためのもので、例えばメモリなどからなる内部記憶装置により実現される。文書記憶部11は、検索対象となる複数の文書を記憶しておくためのもので、例えばハードディスク装置や光ディスク装置などからなる大容量の外部記憶装置により実現される。

【0020】単語辞書12は、多数の単語を、その単語（を構成する文字列）の表記（見出し）、品詞等と対応付けで登録しておくもので、ROM等の不揮発性記憶装置に置かれる。この単語辞書12内の各単語には、その単語の重要度（を示す情報）が付されている。本実施形態において、単語辞書12内の重要度を含む登録情報は固定であるものとする。このような状態で、上記単語の重要度を変更可能とするためには、例えば単語辞書12に優先して使用されるユーザ辞書を設け、このユーザ辞書に所望の単語を登録して、その単語の重要度を任意に設定可能とすればよい。

【0021】単語分割部4、文節合成部5、文節類似度計算部6、最適組み合わせ算出部7、文重み付け部9及び文書類似度計算部10は、文書検索処理に必要な機能要素であり、それぞれ固有の処理ルーチン（と当該処理ルーチンを実行するCPU）により実現されるものである。

【0022】単語分割部4は、入力部1から検索のキーとなる文書を構成する文をなす文字列が入力された場合に、その文字列からなる入力文を検索入力文として単語単位に分割するものである。ここで文とは、句点で区切られる文字列、箇条書きの文字列等、文書を構成する独立の文字列をいう。

【0023】文節合成部5は、単語分割部4により分割された単語からなる文節を合成するものである。文節類似度計算部6は、文節合成部5により得られた検索入力文の文節と、文書記憶部11に記憶されている文書（検索対象文書）の1文（検索対象文）の文節との類似度計算を行うものである。

【0024】最適組み合わせ算出部7は、文節類似度計算部6での類似度計算結果に基づいて、文全体として最も類似している文節の組み合わせを求めるものである。文重み付け部9は、入力文保存部8に保存された入力文

書の各文について、その文書中で使用されている単語の重要度または頻度（使用頻度）に応じて重み付けを行うものである。

【0025】文書類似度計算部10は、入力文保存部8に保存された入力文書のすべての文（検索入力文）について、文書記憶部11に記憶されている文書（検索対象文書）の各文（検索対象文）との類似度を計算して、順次加算し、入力文書（指定文書）と検索対象文書との文書全体としての類似度を算出するものである。

【0026】次に、図1の構成の動作について、その概要を説明する。まず入力部1から、検索のキーとなる文書（入力文書）を構成する文（検索入力文）をなす文字列が入力されると、当該文字列（入力文）は制御部2を介して単語分割部4に送られる。単語分割部4は、単語辞書12を参照して、この入力文を単語単位に分割する。このとき単語分割部4は、分割した単語単位に、品詞、及び活用形、時制などの（形態素情報と呼ばれる）単語属性情報、及び重要度を付加する。

【0027】単語分割部4が入力文を単語単位に分割した結果は、入力文保存部8に保存されると共に、文節合成部5に送られる。文節合成部5は、単語分割部4により分割された単語から文節を合成する。

【0028】文節合成部5が文節合成した結果は、入力文保存部8に送られ、先に保存された単語分割結果に対応付けて保存される。以上の処理は、入力部1から順次入力される各文について繰り返され、単語ごとに単語属性情報及び重要度が付加され、さらに文節合成された状態で、それぞれ検索のキーとなる入力文書（指定文書）の検索入力文として入力文保存部8に保存される。

【0029】さて、入力文保存部8に検索入力文の群からなる入力文書が保存されている状態で、入力部1から、検索要求のコマンドが入力されると、まず制御部2により文重み付け部9が起動される。すると文重み付け部9は、入力文保存部8に保存された各検索入力文それぞれの重みを、対応する入力文書中で使用されている単語の重要度または頻度をもとに計算し、文書類似度計算部10に送る。

【0030】文書類似度計算部10は、文書記憶部11から検索対象となる文書を先頭から1文書取り出し、さらに、その1文書の先頭から1文（検索対象文）を取り出して文節類似度計算部6に送る。また文書類似度計算部10は、入力文保存部8に保存された入力文書中の文（検索入力文）を先頭から1文取り出して文節類似度計算部6に送る。

【0031】文節類似度計算部6は、文書類似度計算部10から送られた入力文の（文書構成単位としての）文節と、検索対象文の（文書構成単位としての）文節とを比較して類似度を計算する動作を、（入力文の文節数）×（検索対象文の文節数）の総当たりで行う。最適組み合わせ算出部7は、この文節類似度計算部6での各文節

ごとの類似度計算結果に基づいて、文全体として最も類似している文節の組み合わせを求める。

【0032】文書類似度計算部10は、最適組み合わせ算出部7での文節組み合わせ結果から文全体の類似度を計算し、その計算結果（類似度）と、文重み付け部9により算出され（て当該文重み付け部9から送られ）た対応する入力文の重みとの積（乗算値）を求め、文書全体の類似度（初期値は0）に加算する。

【0033】文書類似度計算部10は、文書記憶部11から取り出した文書（検索対象文書）に未処理の文（検索対象文）が残っているか否かを調べ、残っているならば次の文を取り出して、上記入力文（検索入力文）との類似度を計算し、その計算結果と対応する入力文の重みとの積を求めて、文書全体の類似度に加算する動作を繰り返す。

【0034】一方、文書記憶部11から取り出した文書に未処理の文が残っていないならば、文書類似度計算部10は、入力文保存部8に次の入力文（未処理の入力文）があるか否かを調べ、あるならば文書記憶部11から取り出した文書の先頭から1文を取り出して、当該次の入力文との類似度を計算し、その計算結果と対応する入力文の重みとの積を求めて、文書全体の類似度に加算する動作を繰り返す。

【0035】やがて、入力文保存部8に次の入力文がなくなると、文書類似度計算部10は、入力文保存部8に保存されている文書（入力文書）と文書記憶部11から取り出した文書（検索対象文書）との間の文書全体の類似度が求まったものと判断する。

【0036】この場合、文書類似度計算部10は、求まった文書全体の類似度から、該当する検索対象文書が入力文書に類似しているか否かを判断し、類似していると判断できたならば、当該検索対象文書、即ち文書記憶部11から取り出した文書を、入力文書に類似した文書であるとして表示部3に表示する。

【0037】次に文書類似度計算部10は、文書記憶部11に未処理の文書（検索対象文書）が残っているか否かを調べ、残っているならば次の文書を取り出す。そして文書類似度計算部10は、入力文保存部8に保存された入力文書中の文（検索入力文）を先頭から1文取り出すと共に、上記文書記憶部11から取り出した検索対象文書の先頭から1文（検索対象文）を取り出して、入力文と検索対象文との類似度を求め、その類似度と対応する入力文の重みとの積を求めて、その検索対象文書（と入力文書）との間の文書全体の類似度に加算する動作を繰り返す。

【0038】やがて、文書記憶部11内に未処理の文書がなくなると、上記した一連の検索処理は終了する。以上が、図1の構成の動作の概要である。

【0039】次に、図1の構成の動作の詳細について、図2乃至図4のフローチャートを参照して説明する。ま

ず制御部2は、入力部1から検索キーとなる文書の1文（検索入力文）をなす文字列が入力されたと判断すると（ステップS1）、当該入力文（入力文字列）を入力し（ステップS2）、表示部3の表示画面に表示すると共に、単語分割部4に送る。

【0040】ここでは、入力文として、図5（a）において符号a1で示すような「A社は、低価格のノート型EWSを開発した。」という文字列が入力され、単語分割部4に送られたとする。

【0041】単語分割部4は、制御部2から送られた入力文（検索入力文）a1の文字列を単語辞書12の登録情報を用いて解析することで単語単位に分割する（ステップS3）。

【0042】ここでは単語分割部4は、図5（a）の入力文a1を、図6（a）に示すように、「A社」「は」「低価格」「の」「ノート型」「EWS」「を」「開発する」のように単語に分割し、それぞれの品詞、活用形及び時制などの単語属性の情報、及び（単語辞書12に登録されている）当該単語の重要度を付加する。この図6（a）に示したような単語分割部4による単語分割結果は、入力文保存部8に保存されると共に文節合成部5に渡される。

【0043】単語分割部4による単語分割処理が終了すると、制御部2により文節合成部5が起動される。文節合成部5は、単語分割部4により得られた単語を合成して、文節を生成する（ステップS4）。

【0044】ここでは、図6（a）に示す単語分割処理結果から、図6（b）のように、「A社は」「低価格の」「ノート型」「EWSを」「開発する」の5文節が生成される。

【0045】文節合成部5による文節単位への合成処理が終了すると、図6（b）に示したような当該文節合成の結果が、図6（a）に示したような単語分割結果に対応付けて、制御部2により入力文保存部8に保存される（ステップS5）。

【0046】以上のステップS1からS5までの処理を繰り返すことで、入力文保存部8には入力部1から入力された入力文（検索入力文）が図6（a）、（b）に示したような形式で順次保存される。

【0047】図5（a）は、このような入力部1から順次文を入力してできた入力文書を示すもので、上記した入力文a1の他に、入力文a2、a3から構成されるものとする。

【0048】次に、入力文保存部8に入力文書が保存されている状態で、入力部1から検索要求のコマンドが入力されたものとする。制御部2は、入力部1から検索要求のコマンドが入力されたことを検出すると（ステップS1）、文重み付け部9を起動する。すると文重み付け部9は、入力文保存部8に保存されている入力文書の各文（検索入力文）それぞれの重みを、その文書中で使用

されている単語の重要度または頻度をもとに算出する重み付け処理を行う（ステップS6）。

【0049】以下、この文重み付け部9によるステップS6の重み付け処理の詳細を、入力文保存部8に保存されている文書中で使用されている単語の重要度から各文に重み付けをする場合を例に、図4（a）のフローチャートを参照して説明する。

【0050】文重み付け部9は、入力文保存部8に保存されている入力文書の各文ごとに、その文で使われている単語の重要度を加算して、各文ごとの重みを求める（ステップS21）。

ここで各文で使われている単語の重要度は、前記ステップS3で単語分割部4が単語単位の分割を行った時点で、図6（a）の例のように単語ごとに付加されて、入力文保存部8に保存されたものである。

【0051】図5（a）に示した入力文a1～a3からなる文書が入力文保存部8に保存されている例では、先頭の「A社は、低価格のノート型EWSを開発した。」という入力文a1中の各単語の重要度は、図6（a）から明らかなように、「A社」=1、「ノート型」=2、「EWS」=1なので（他の単語の重要度はすべて0）、その重要度を加算した結果は $1+2+1=4$ となり、これが先頭の入力文a1の重みとされる。

【0052】同様にして、後続の入力文a2、a3についても、その入力文a2、a3中の各単語の重要度を加算する処理が行われて、その加算結果が入力文a2、a3の重みとされる。

【0053】図7（a）は、このようにして求められた（入力文保存部8に保存されている文書中の）各入力文a1～a3の重みの一例を示す。次に、上記ステップS6の処理の詳細を、入力文保存部8に保存されている文書に使われている単語の頻度から各文に重み付けをする場合を例に、図4（b）のフローチャートを参照して説明する。

【0054】まず文重み付け部9は、入力文保存部8に保存されている文書（ここでは、図5（a）に示すように文a1～a3からなる文書）で使用されている各単語の頻度（使用頻度）を調べる（ステップS31）。

【0055】ここでは、各単語の頻度を上記した単語の重要度と同様に扱うようにしており、1文書中で1回しか使用されない単語（頻度1の単語）については対象外とする。したがって、図5（a）の文書の例では、各単語の頻度として、「A社」=2、「ノート型」=2、「EWS」=3、「B社」=2が求められる。

【0056】次に文重み付け部9は、入力文保存部8に保存されている文書の各文ごとに、その文で使われている各単語の頻度を加算して、各文ごとの重みを求める（ステップS32）。

【0057】図5（a）に示した入力文a1～a3からなる文書が入力文保存部8に保存されている例では、先

頭の「A社は、低価格のノート型EWSを開発した。」という入力文（第1文）a1の重みは、当該入力文a1中の単語の頻度が「A社」=2、「ノート型」=2、「EWS」=3であることから、 $2+2+3=7$ となる。

【0058】同様に、「ノート型EWSは、すでにB社が今年5月に発売した。」という入力文（第2文）a2の重みは、当該入力文a2中の単語の頻度が「ノート型」=2、「EWS」=3、「B社」=2であることから、 $2+3+2=7$ となり、「A社のEWSは、B社よりも10%程度価格が下がった。」という入力文（第3文）a3の重みは、当該入力文a3中の単語の頻度が、「A社」=2、「EWS」=3、「B社」=2の頻度から $2+3+2=7$ となる以上が文重み付け部9によるステップS6の処理の詳細である。

【0059】文重み付け部9によるステップ6の処理が行われて、入力文保存部8に保存されている文書の各文ごとの重みが求められると、制御部2により文書類似度計算部10が起動される。

【0060】文書類似度計算部10は、文書記憶部11から検索対象となる文書を先頭から順に1文書取り出す（ステップS7）。ここでは、文書記憶部11から、図5（b）に示すような、文b1、b2からなる先頭の文書が取り出されたものとする。但し、文書記憶部11に実際に格納されている文書（中の各文）の形式は、入力文保存部8に保存されている文書（中の各文）と同様であり、図5（b）に示したような単なる文字列（文字コード列）の形ではない。即ち文書記憶部11には、各文書中の各文が、あらかじめ単語分割部4の処理（または単語分割部4と同等の機能）により分割された（単語属性情報及び重要度が付加された）単語列と、文節合成部5の処理（または文節合成部5と同等の機能）により生成された文節列の形で格納されている。

【0061】文書類似度計算部10は、ステップS7で文書記憶部11から検索対象となる1文書を取り出すと、今度は入力文保存部8から先頭の1文（検索入力文）を取り出すと共に（ステップS8）、ステップS7で文書記憶部11から取り出した検索対象文書の先頭から1文を検索対象文として取り出し（ステップS9）、その取り出した検索入力文と検索対象文の文節を文節類似度計算部6に渡す。

【0062】文節類似度計算部6は、文書類似度計算部10がステップS8で取り出した検索入力文の文節とステップS9で取り出した文書記憶部11中の文書（検索対象文書）の1文（検索対象文）の文節とを当該文書類似度計算部10から受け取ると、その検索入力文の1文節と検索対象文の1文節とを比較して、その文節間の類似度を計算する（ステップS10）。この文節類似度計算部6での類似度計算は、（検索入力文の文節数）×（検索対象文の文節数）の総当たりで行われ、すべての

文節について繰り返される（ステップS11）。

【0063】最適組み合わせ算出部7は、文節類似度計算部6による検索入力文と検索対象文との文節間の類似度計算の結果に基づいて、当該検索入力文と検索対象文との間で文全体として最も類似している文節の組み合わせ（最適組み合わせ）を求める（ステップS12）。

【0064】ここでは、図5（a）に示す文書中の入力文（検索入力文）a1と、図5（b）に示す文書中の第1文（検索対象文）b1の文節間の類似度計算が行われ、その文a1、b1間で文全体として最も類似している文節の組み合わせとして、「A社は」と「A社は」、「低価格の」と「低価格の」、及び「EWSを」と「EWSを」の3文節が求められたものとする。この最適組み合わせ算出部7での最適組み合わせ算出結果は、文書類似度計算部10に渡される。

【0065】文書類似度計算部10は、最適組み合わせ算出部7から最適組み合わせ算出結果を受け取ると、その結果をもとに、例えば検索入力文中の文節数に対する最も類似している文節の組み合わせ数の割合を求めて文全体の類似度とし、それにステップS6で文重み付け部9により求められている対応する検索入力文の重みを乗じた値（文全体の類似度と検索対象文書の重みとの積）を算出する（ステップS13）。

【0066】ここでは、最適組み合わせ算出部7の処理より、検索入力文a1の5文節中「A社は」「低価格の」「EWSを」の3文節が検索対象文b1と一致していることが求められているため、その割合（最も類似している文節の組み合わせ数の割合） $3/5=0.6$ が、検索入力文a1と検索対象文b1との間の文全体の類似度として計算される。また、文重み付け部9で求められた検索入力文a1の重みは、例えば単語の重要度を用いた重み付けの例では、4である。この場合、ステップS13で求められる、検索入力文a1と検索対象文b1の間の類似度と検索入力文a1の重みとの積は、 $0.6 \times 4 = 2.4$ となる。

【0067】次に文書類似度計算部10は、ステップS13で求めた値（文全体の類似度と重みとの乗算値）を、文書全体の類似度（ここでは、入力文保存部8に保存されている文書とステップS7で文書記憶部11から取り出された検索対象文書との間の文書全体の類似度）に加算する（ステップS14）。この文書全体の類似度の初期値は0であり、したがって1回目の加算結果は、ステップS13で求めた値に一致する。

【0068】文書類似度計算部10はステップS14を実行すると、ステップS7で取り出した検索対象文書の中に未処理の文、即ちステップS9で未だ取り出されていない文が残っているか否かを調べる（ステップS15）。

【0069】もし、未処理の文があるならば、文書類似度計算部10はステップS9に戻り、ステップS7で取

り出した検索対象文書の中から次の1文を検索対象文として取り出す。ここで取り出される次の文は、図5

(b)に示す文書中の第2文(最後の検索対象文)b2である。

【0070】さて、検索入力文a1と検索対象文b2の間では、「ノート型の」と「ノート型の」、「EWS」と「EWS」の2文節が対応している。このため、上記ステップS9に続くステップS10～S12を経た後、検索入力文a1中の文節数に対する最も類似している文節の組み合わせ数の割合 $2/5=0.4$ が文全体の類似度として求められ、これに検索入力文a1の重み4を乗じた値 $0.4 \times 4 = 1.6$ が(その時点における)文書全体の類似度に加算される(ステップS13, S14)。ここでは、加算前の文書全体の類似度が2.4であることから、新たな文書全体の類似度は $2.4 + 1.6 = 4.0$ となる。

【0071】ステップS14が実行されると、前記したように、ステップS7で取り出した検索対象文書の中に未処理の文が残っているか否かが調べられる(ステップS15)。

【0072】この例のように未処理の文がないならば、文書類似度計算部10は、入力文保存部8に保存されている文書の中に未処理の文、即ちステップS8で取り出した検索入力文が残っているか否かを調べる(ステップS16)。

【0073】もし、未処理の文があるならば、文書類似度計算部10はステップS8に戻り、入力文保存部8に保存されている文書の中から次の1文を検索入力文として取り出す。ここで取り出される次の文は、図5(a)に示す文書中の第2文(検索入力文)a2である。

【0074】文書類似度計算部10は、ステップS8で次の検索入力文a2を取り出すと、前記ステップS9以降の処理を繰り返す。これにより、図5(a)中の検索入力文a2に対して図5(b)中の各検索対象文b1, b2の類似度が順次計算され、その都度当該類似度に検索入力文a2の重み(ここでは図7(a)に示すように3)を乗じた値が求められて、文書全体の類似度に加算される。

【0075】以下、同様にして、入力文保存部8に保存されている文書中の第3文(最後の検索入力文)a3がステップS8で取り出され、ステップS9以降の処理が繰り返される。すると、図5(a)中の検索入力文a3に対して図5(b)中の各検索対象文b1, b2の類似度が順次計算され、その都度当該類似度に検索入力文a3の重み(ここでは図7(a)に示すように2)を乗じた値が求められて、文書全体の類似度に加算される。

【0076】この例では、検索入力文a1～a3と検索対象文b1, b2との間の文全体の類似度(検索入力文中の文節数に対する、検索入力文と検索対象文との間で最も類似している文節の組み合わせ数の割合)は図7

(b)のようになる。

【0077】そして、検索入力文と検索対象文との間の文全体の類似度が求められる都度、上記ステップS13, S14にて、その類似度と当該検索入力文の重みとの積の値が、文書全体の類似度に加算されていく。

【0078】したがって、図7(b)に示したような検索入力文a1～a3と検索対象文b1, b2との間の文全体の各類似度に、対応する検索入力文の重み(図7(a)参照)を乗じ、その乗算値(類似度と重みとの積)を順次加算し終えた段階では、第5(a)に示した入力文書と図5(b)に示した検索対象文書との間の文書全体の類似度は、

$0.6 \times 4 + 0.4 \times 4 + 0.25 \times 3 + 0.375 \times 3 + 0.14 \times 2 + 0.57 \times 2 = 7.295$ のように求められる。

【0079】この段階では、ステップS7で取り出した検索対象文書の中に未処理の文はなく(ステップS15)、また入力文保存部8に保存された文書の中にも未処理の文はないことから(ステップS16)、文書類似度計算部10はステップS17に進む。

【0080】文書類似度計算部10は、ステップ17において、その時点において求められている文書全体の類似度(ここでは、7.295)から、入力文保存部8に保存された入力文書と文書記憶部11から取り出された検索対象文書とが類似しているか否かを判断する。ここでの判断基準は、設定された基準値(閾値)を越える類似度を持つか否かであり、この基準値は、ユーザの要求する検索精度に応じて適宜変更設定可能である。例えば、最低の検索精度でよい場合には、上記基準値は0に設定される。この場合、文書全体の類似度として0以外の値が求められているならば、即ち入力文書と検索対象文書との間に少しでも類似度を持つならば、その大小に無関係に両文書は類似していると判断される。

【0081】文書類似度計算部10は、入力文保存部8に保存された文書と文書記憶部11から取り出された検索対象文書とが類似していると判断した場合、その旨を制御部2に通知する。すると制御部2は、文書類似度計算部10により入力文書に類似していると判断された検索対象文書を表示部3に検索結果として表示する(ステップS18)。

【0082】次に文書類似度計算部10は、文書記憶部11内に未処理の文書、即ちステップS7で未だ取り出されていない文書が残っているか否かを調べる(ステップS19)。

【0083】もし、未処理の文書があるならば、文書類似度計算部10はステップS7に戻り、文書記憶部11から次の1文書を検索対象文書として取り出し、前記したステップS8以降の処理を行う。

【0084】これに対し、未処理の文書がないならば、文書類似度計算部10は制御部2に制御を戻して処理を

終了する。なお、本発明は前記した実施の形態に限定されるものではない。

【0085】即ち、前記実施形態では、入力文保存部8に保存された入力文書の各文の重み付けを行うのに、単語の重要度に従う重み付けの手法を適用した場合と、入力文書中で使用されている単語の頻度（使用頻度）に従う重み付けの手法を適用した場合について説明したが、これに限るものではない。例えば、単語の重要度をすべて同じ値（例えば1）に設定すれば、重み付けしない場合と等価であると見なせる。また、単語の重要度と単語の頻度とを組み合わせる重み付けするようにしても構わない。具体的には、入力文書に現れる単語の重要度と当該単語の頻度との積を各文で加算し、その加算結果を該当する文の重みとするようにしてもよい。この場合、頻度1の単語も対象とするとよい。

【0086】また、前記実施形態では、入力文保存部8に保存される（単語分割、文節合成後の）入力文の群からなる文書、即ち検索キーとなる文書が、入力部1から入力されたものである場合について説明したが、例えば文書記憶部11に記憶されている複数の文書の中からユーザにより指定された文書であっても構わない。

【0087】また、前記実施形態では、文書記憶部11に記憶されている複数の文書のすべてが検索の範囲（検索対象文書）となる場合について説明したが、ユーザにより指定された範囲の文書だけを検索対象文書とすることも可能である。要するに本発明は、その要旨を逸脱しない範囲で種々変形して実施することができる。

【0088】

【発明の効果】以上詳述したように本発明によれば、指定文書（ユーザの入力した文書あるいはユーザの指定した文書）中のすべての文を検索入力文として、検索対象として用意された文書中の各文に対して、1文全体の類似度を求めて順次加算した値を文書全体の類似度とし、指定文書と検索対象文書とが文書全体としてどれだけ類似しているかを、当該文書全体の類似度で表すことにより、この文書全体の類似度をもとに検索対象となる複数の文書の中から指定文書と類似した文書を検索することができ、ユーザが望む文書と類似する文書を検索する際に、ユーザは検索キーとして単語やキーワードなどを選び出すというような手間が一切必要なくなる。

【0089】また本発明によれば、指定文書で使用されている各単語の重要度または使用頻度の少なくとも一方に応じて当該指定文書中の各文に重み付けを行い、その各文ごとの重みをその文（検索対象文）と検索入力文との文全体の類似度に反映させることで、すべての検索対象文書の中から指定文書と類似した文書であって且つ指

定文書中で使用している重要度の高い単語を多数含む、あるいは指定文書中で使用している単語の使用数の多い文書を検索することができる。

【図面の簡単な説明】

【図1】本発明の一実施形態に係る文書検索装置の全体構成を示すブロック図。

【図2】同実施形態の動作を説明するためのフローチャートの一部を示す図。

【図3】同実施形態の動作を説明するためのフローチャートの残りを示す図。

【図4】図3中のステップ6の処理（入力文の重み付け処理）の詳細を説明するためのフローチャートであり、図4（a）は単語の重要度により文に重み付けをする場合のフローチャート、図4（b）は入力文書内での単語の使用頻度により文に重み付けをする場合のフローチャート。

【図5】入力部1から入力された検索のキーとなる入力文書を構成する複数の入力文の例と、文書記憶部11から取り出された検索対象文書を構成する複数の文（検索対象文）の例を示す図。

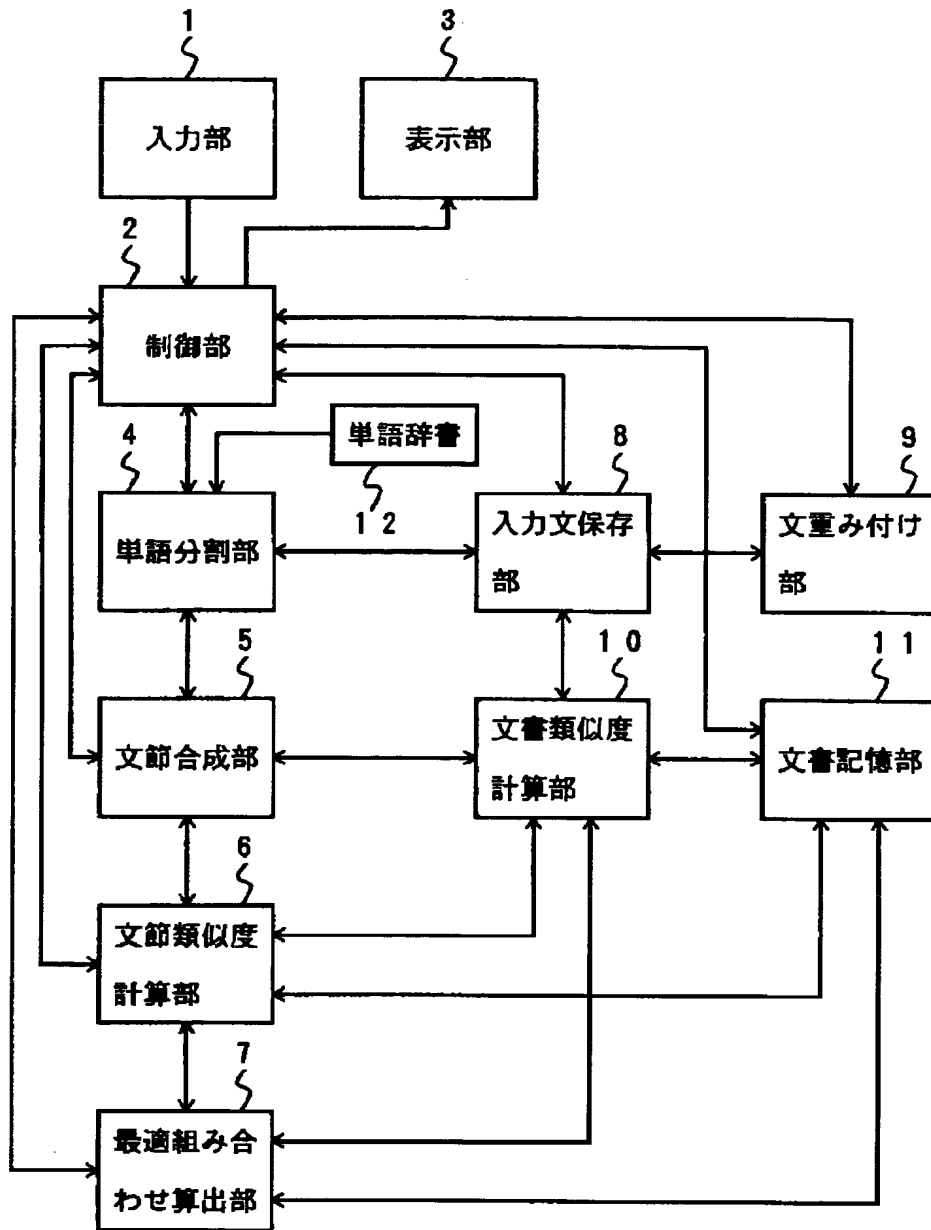
【図6】図1中の単語分割部4による単語分割結果（単語分割された文、当該文の各単語、及び重要度を含む付加情報）の例と、この分割された各単語が文節合成部5により合成されて文節が生成された例を示す図。

【図7】図5中の入力文書の各入力文a1～a3に文重み付け部9が単語の重要度により重み付けした場合の各入力文a1～a3の重みの例と、文書類似度計算部10により求められる図5中の入力文書の各入力文a1～a3と検索対象文b1、b2との間の文全体の各類似度の例を示す図。

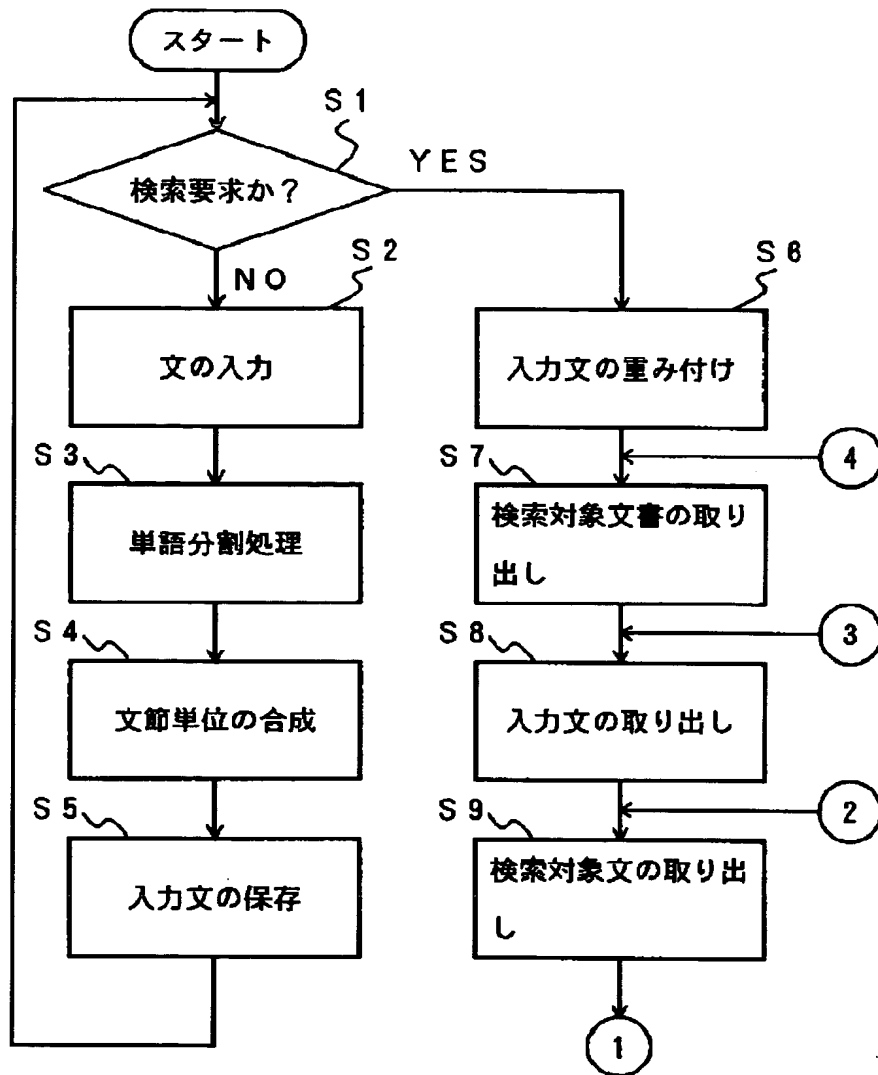
【符号の説明】

- 1…入力部、
- 2…制御部、
- 3…表示部、
- 4…単語分割部、
- 5…文節合成部、
- 6…文節類似度計算部（文構成単位類似度計算手段）、
- 7…最適組み合わせ算出部、
- 8…入力文保存部、
- 9…文重み付け部、
- 10…文書類似度計算部、
- 11…文書記憶部、
- 12…単語辞書、
- a1～a3…文（検索入力文）、
- b1、b2…文（検索対象文）。

【図 1】



【図2】



【図7】

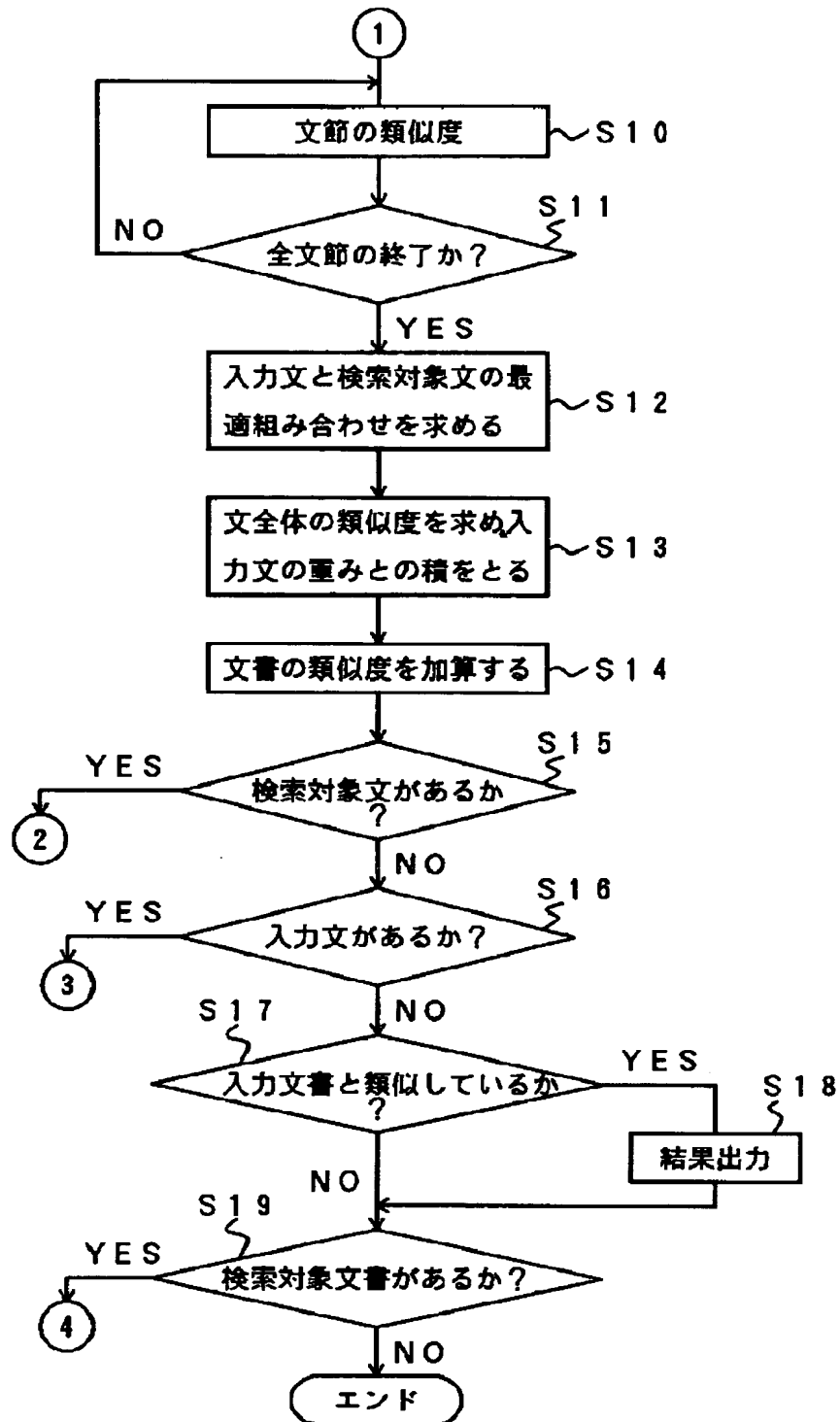
(a)

入力文	重み
a 1	4
a 2	3
a 3	2

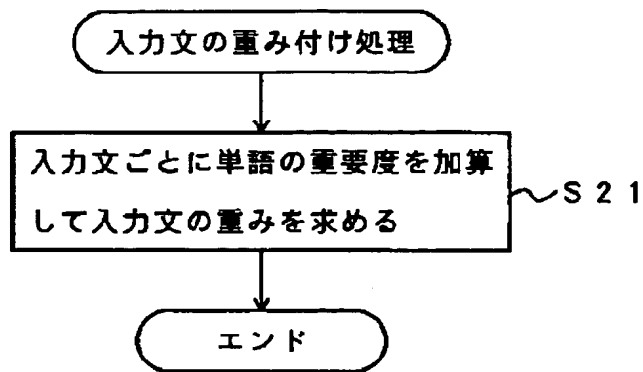
(b)

入力文 \ 検索対象文	b 1	b 2
a 1	0.6	0.4
a 2	0.25	0.375
a 3	0.14	0.57

【図3】



【図4】



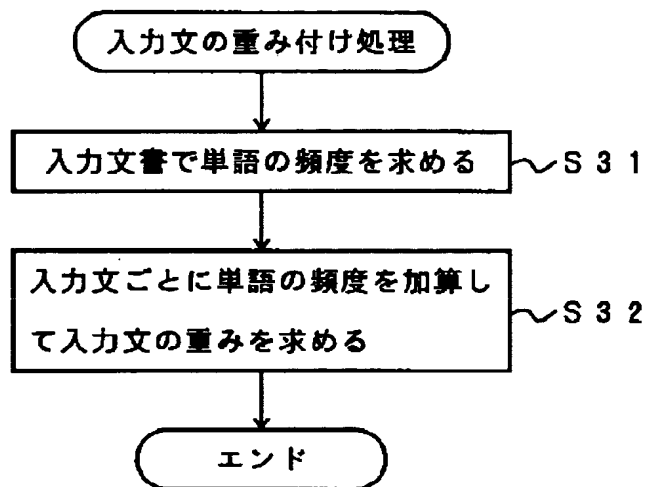
(a)

【図6】

A社は低価格のノート型のEWSを開発した。

No.	単語	品詞	重要度
1	A社	名詞	1
2	は	助詞	0
3	低価格	形容動詞	0
4	の	助詞	0
5	ノート型	名詞	2
6	の	助詞	0
7	EWS	名詞	1
8	を	助詞	0
9	開発する	動詞	0

(a)

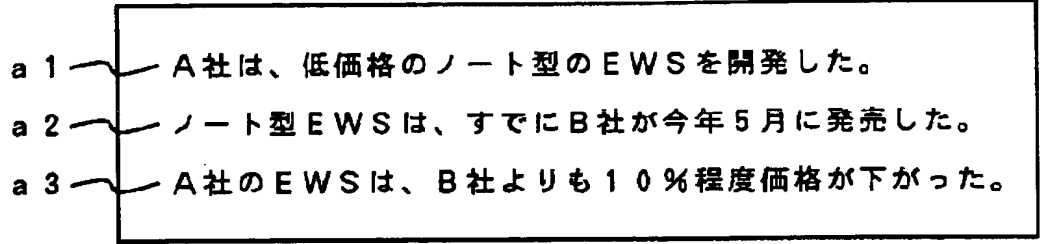


(b)

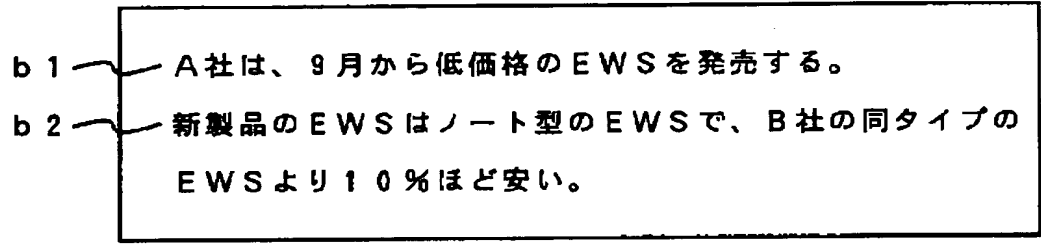
1	A社は
2	低価格の
3	ノート型の
4	EWSを
5	開発する

(b)

【図5】

- 
- a 1 — A社は、低価格のノート型のEWSを開発した。
- a 2 — ノート型EWSは、すでにB社が今年5月に発売した。
- a 3 — A社のEWSは、B社よりも10%程度価格が下がった。

(a)

- 
- b 1 — A社は、9月から低価格のEWSを発売する。
- b 2 — 新製品のEWSはノート型のEWSで、B社の同タイプのEWSより10%ほど安い。

(b)